

# Atelier Data Science

---

Deep learning practice 2  
Convolutional Neural Networks

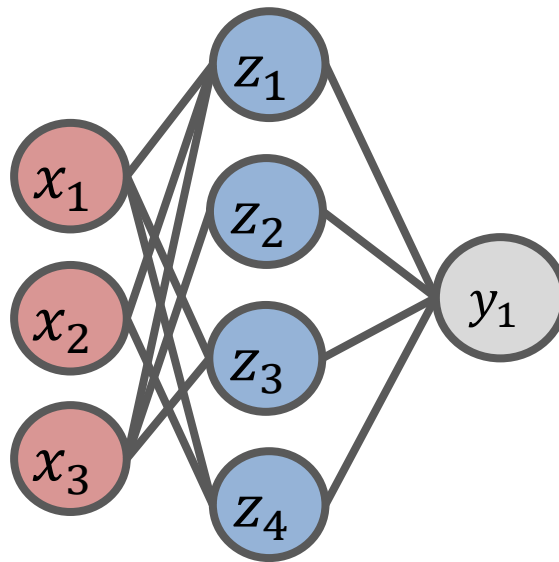
Irina Proskurina

[Irina.Proskurina@univ-lyon2.fr](mailto:Irina.Proskurina@univ-lyon2.fr)

Laboratoire ERIC – Université Lyon 2

# Previous Lesson Recap 1

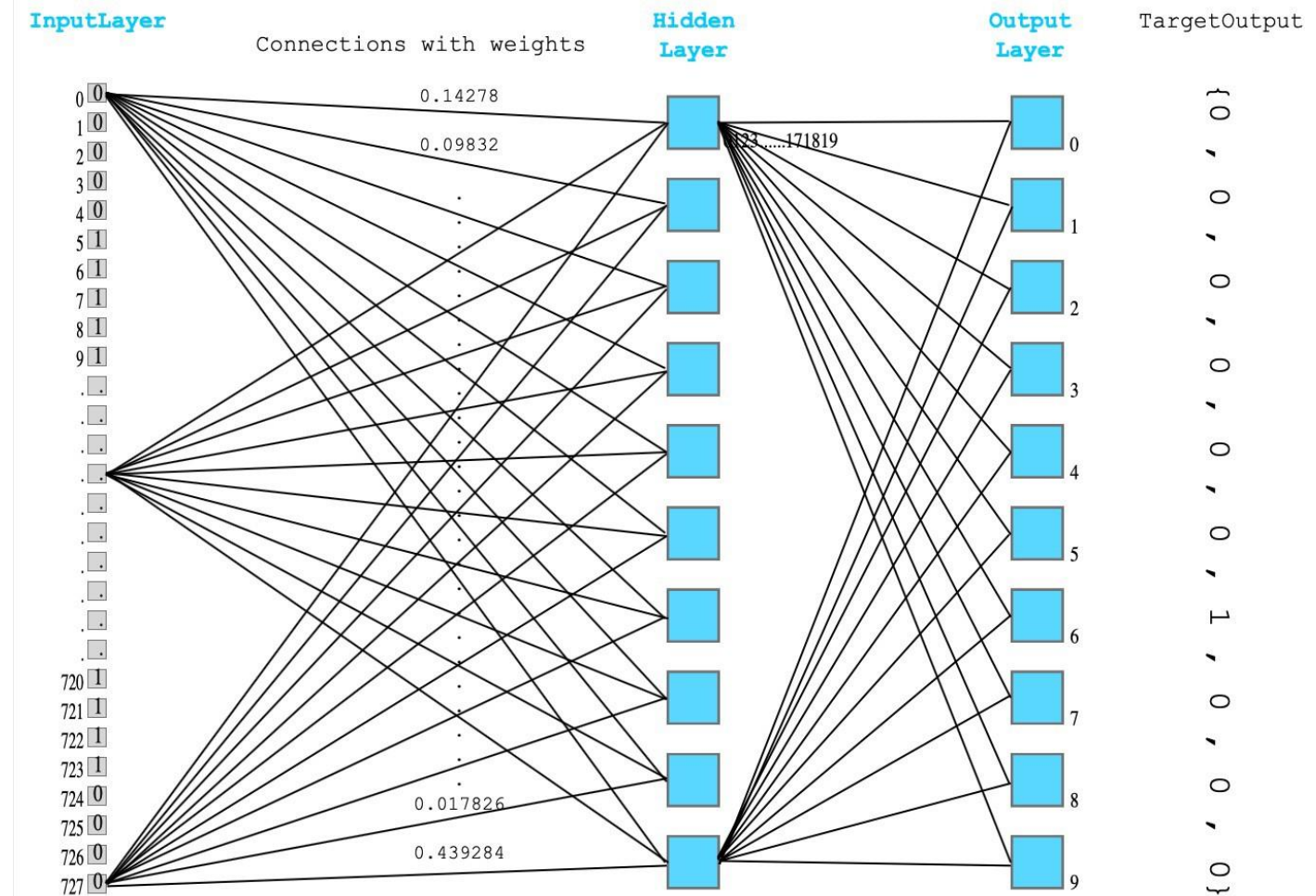
- 1) Neural Networks -> no manual feature engineering is needed!
- 2) Fully connected layers (**torch.nn.Linear**(20, 30))
- 3) Fully connected neural networks: connect every neuron in one layer to every neuron in the other layer + activation function



# Previous Lesson Recap 2

## MNIST

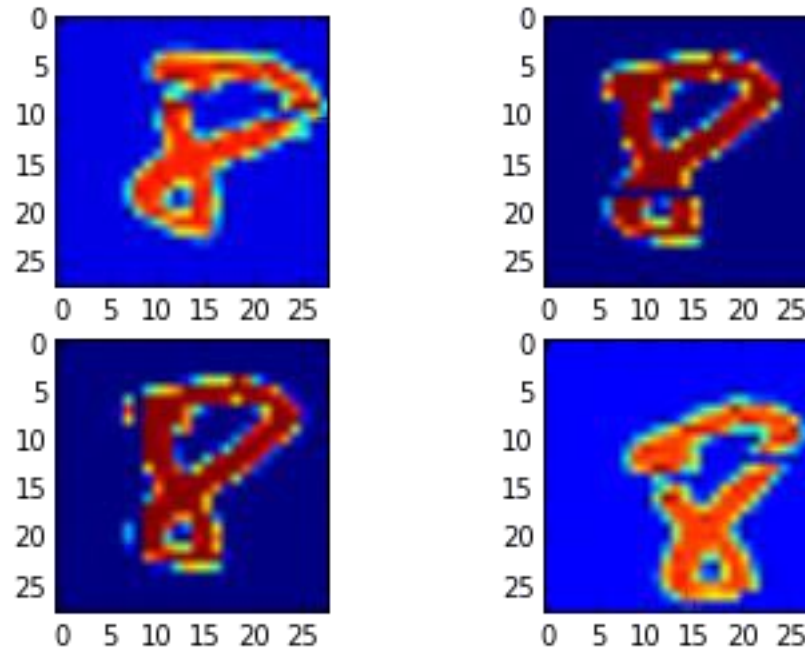
- Each neuron can detect the presence of a specific set of pixels



# Previous Lesson Recap 2

## MNIST

- If you shift the digit slightly, the neuron will no longer detect its pattern



# Number of parameters

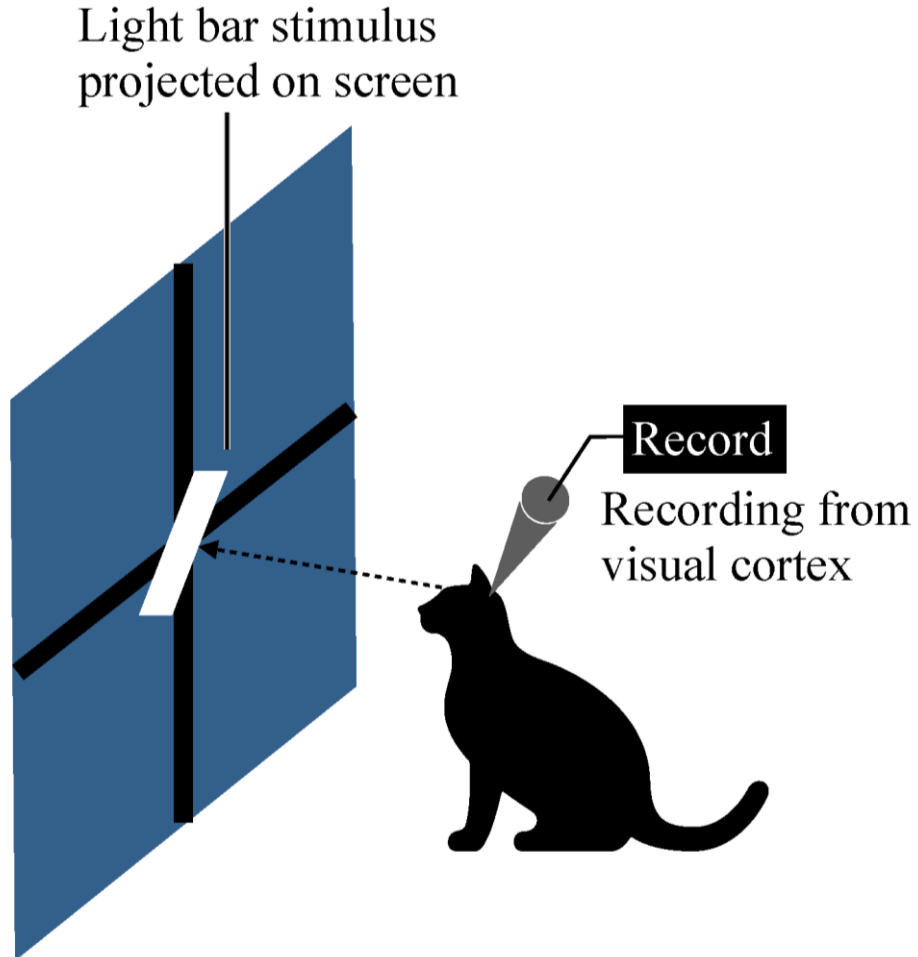
- 784 inputs
- Fully connected layer: 1000 neurons
- Output layer: 10 neurons (one for each class)
- Weights between input and fully connected layers:  
 $(784 + 1) * 1000 = 785,000$
- Weights between fully connected and output layers:  
 $(1000 + 1) * 10 = 10,010$

# Fully connected neural networks for image classification

- A lot of parameters
- Prone to overfitting
- Do not consider special patterns in images: shifts, slight changes in shape, etc.
- One of the best ways to combat overfitting is to reduce the number of parameters

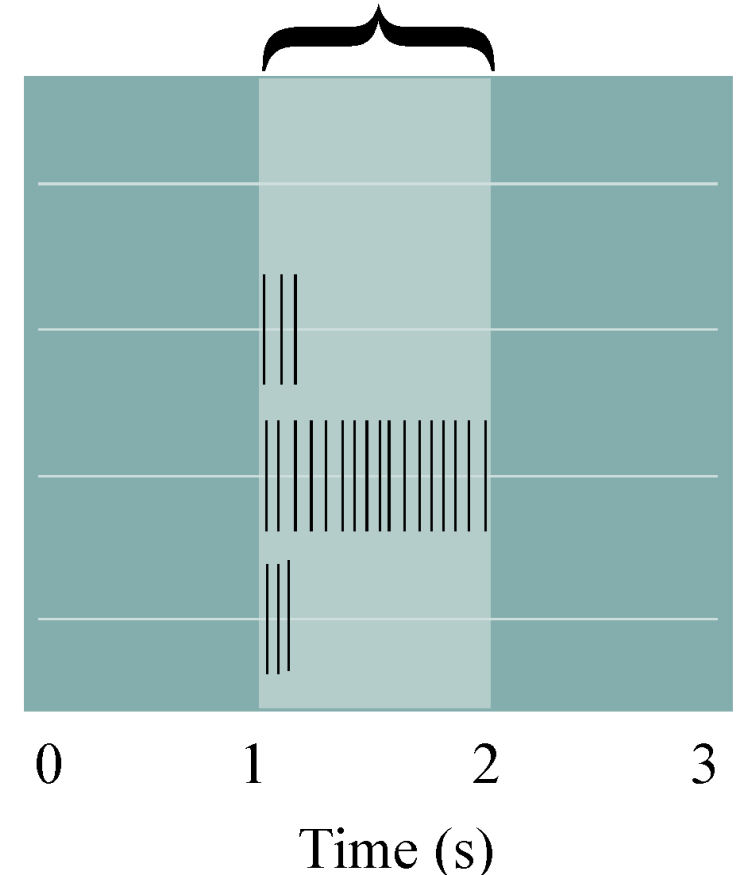
# Convolutional neural network

# Experiments with the visual cortex



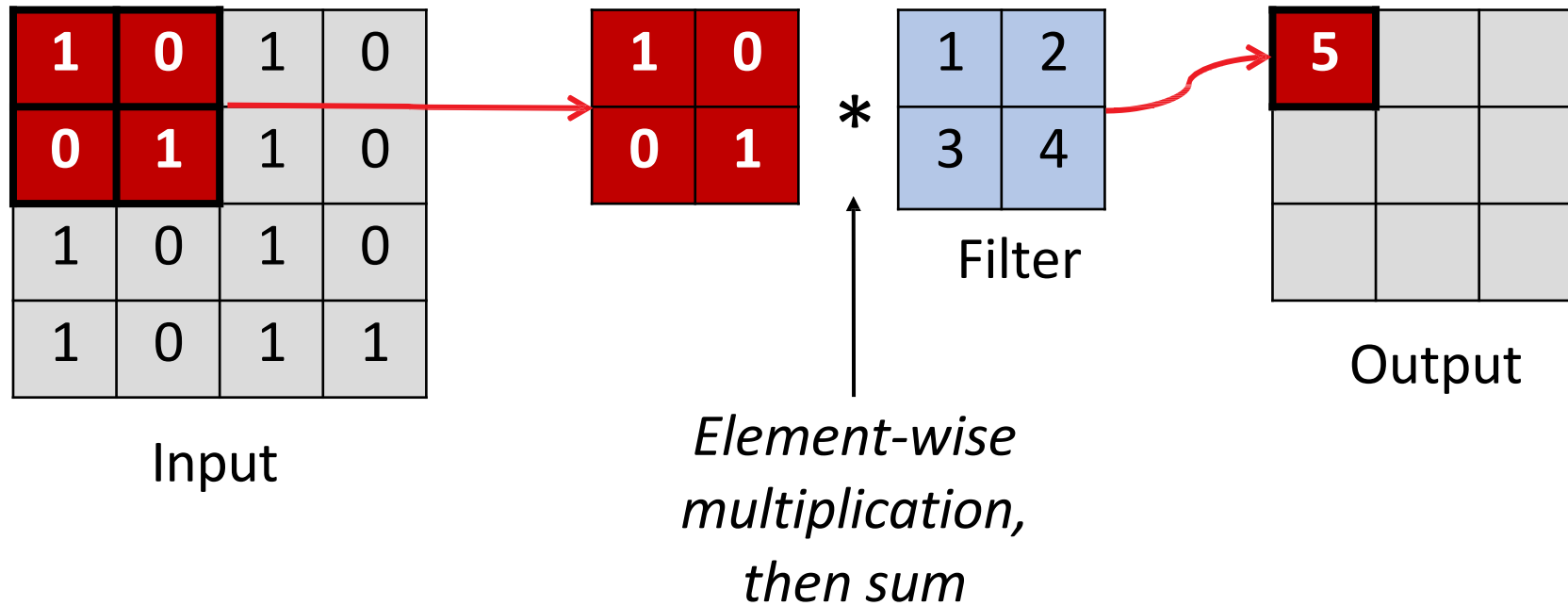
Simulation orientation

Stimulus presented





# Convolution



# Convolution

$$\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{2}$$

$$\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{2}$$

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{1}$$

$$\begin{array}{|c|c|} \hline 0 & 2 \\ \hline 3 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{0}$$

$$\begin{array}{|c|c|} \hline 3 & 0 \\ \hline 0 & 3 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{6}$$

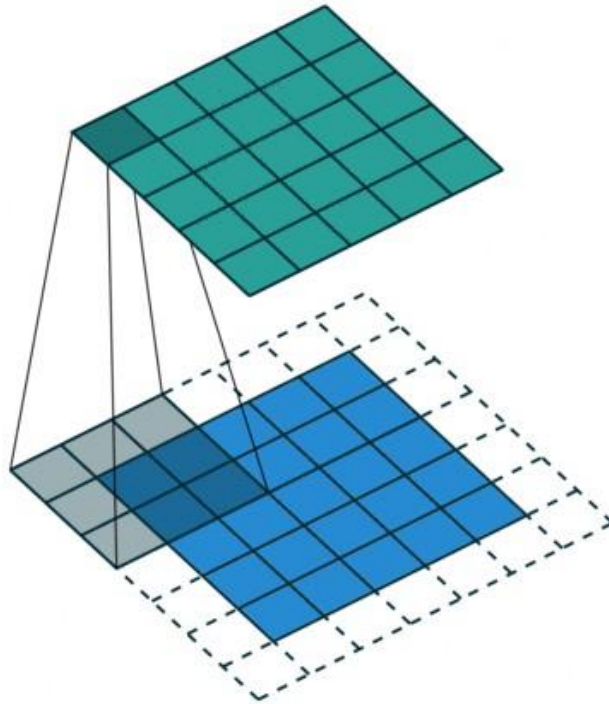
$$\begin{array}{|c|c|} \hline 5 & 0 \\ \hline 0 & 5 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{10}$$

# Convolution

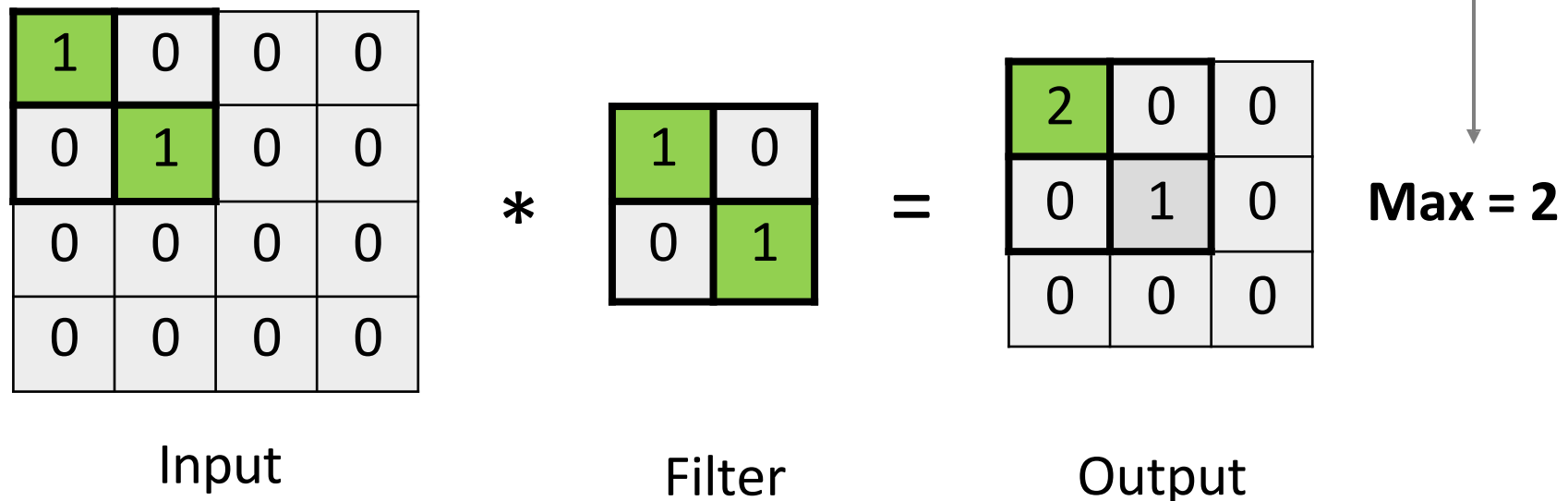
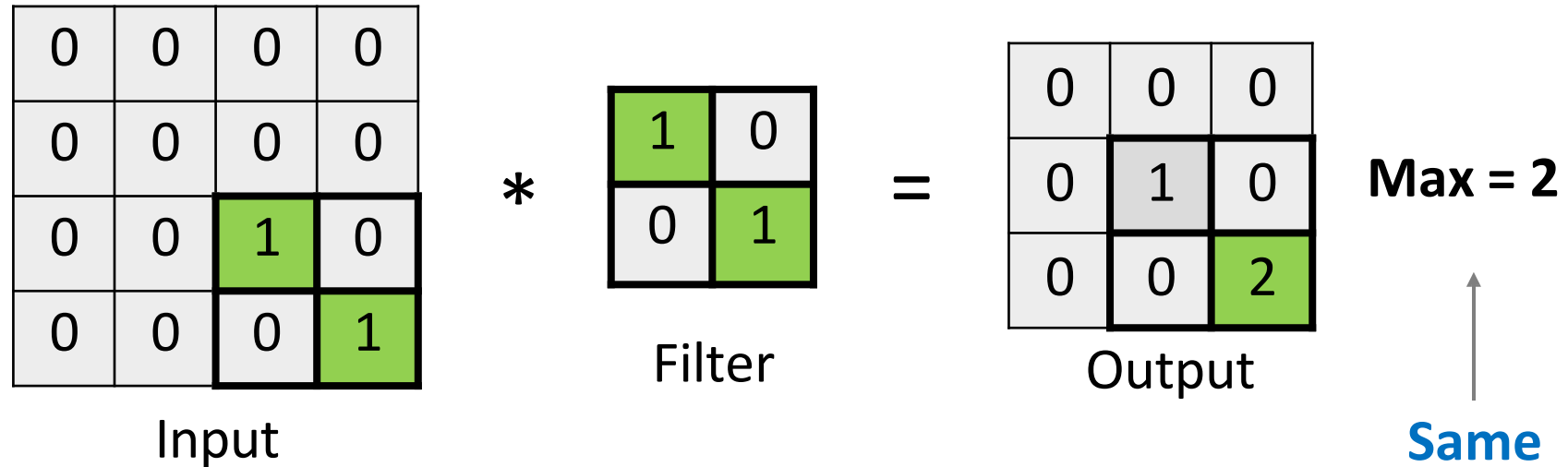
- Detects a pattern in the image, which is defined by a filter
- The stronger the pattern in a particular area of the image, the higher the convolution value will be

# Convolution

- The result of the convolution is a new image, where each pixel is a weighted sum of the pixels in the original image



# The maximum of convolution is invariant to shifts

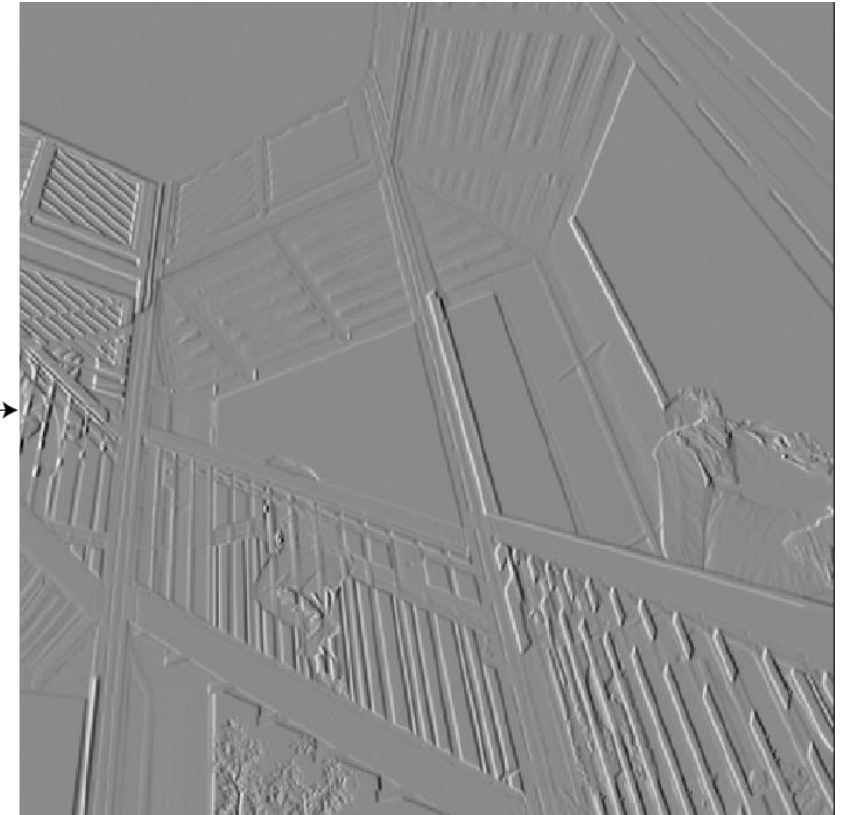


# Convolutions in computer vision



$$\begin{bmatrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{bmatrix}$$

Horizontal Sobel kernel

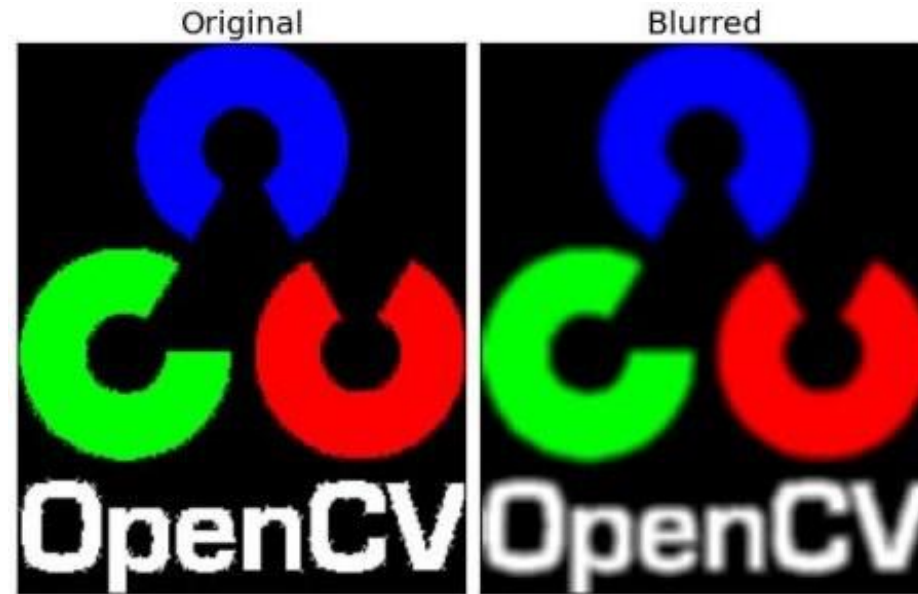


# Convolutions in computer vision



$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} + \begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix} = \begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 2 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

# Convolutions in computer vision



$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$



# Convolution

$$Im^{out}(x, y) = \sum_{i=-d}^d \sum_{j=-d}^d (K(i, j) Im^{in}(x + i, y + j) + b)$$

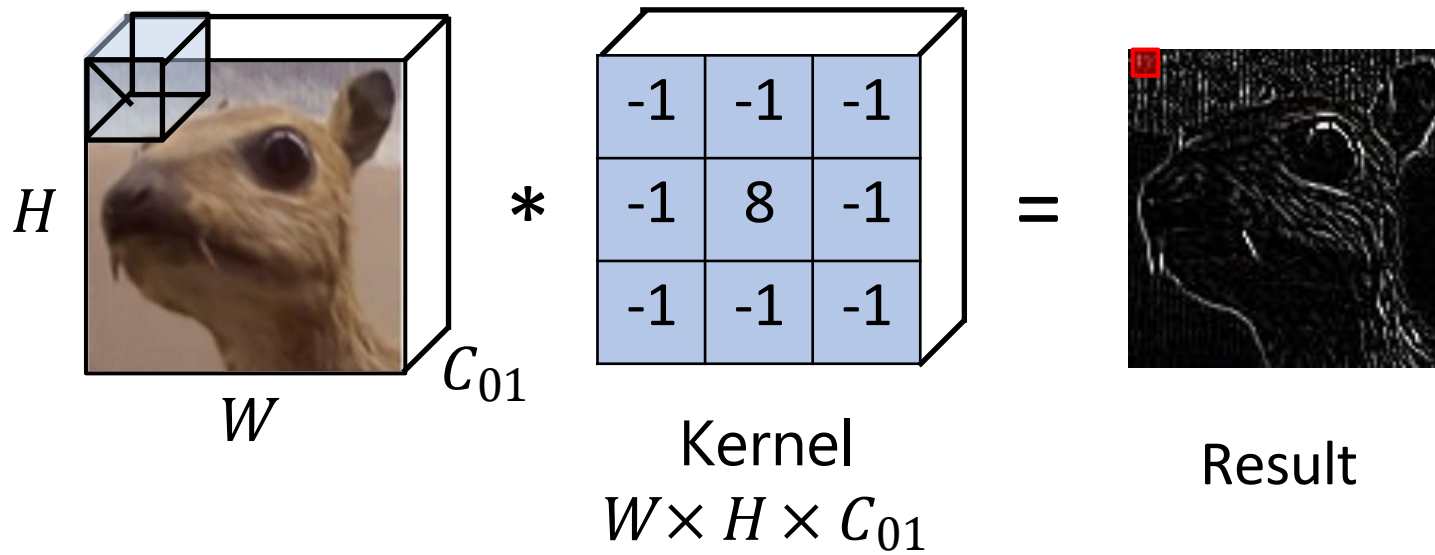
- A pixel in the resulting image depends only on a small set of the input image (local connectivity)
- The weights (kernel values) are the same for all pixels in the output image (shared weights)

# Convolution

- Typically, the original image is in color, implying it has multiple channels (R, G, B)
- Let's take this into account in the formula:

$$Im^{out}(x, y) = \sum_{i=-d}^d \sum_{j=-d}^d \sum_{c=1}^C (K(i, j, c) Im^{in}(x + i, y + j, c) + b)$$

# Convolution



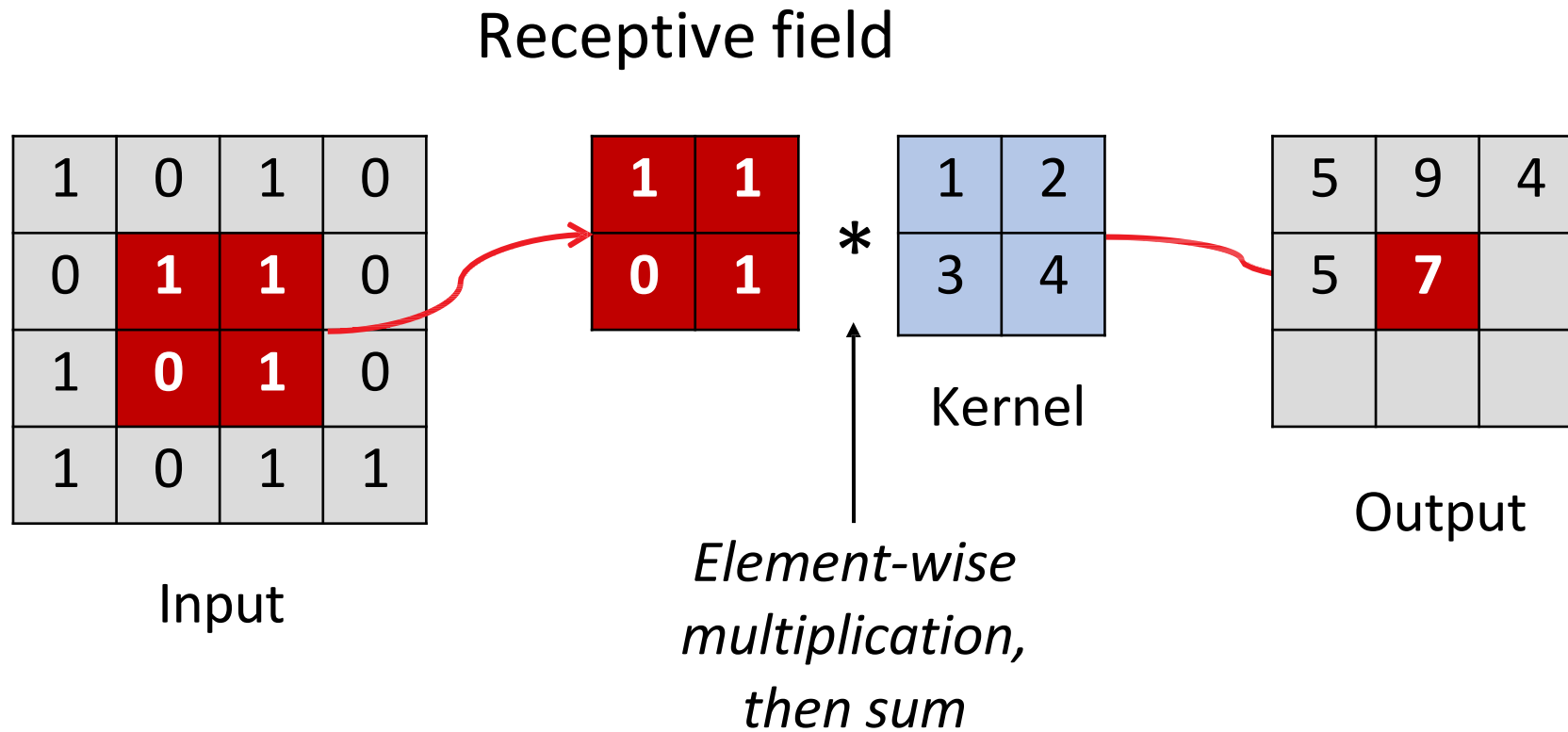
# Number of parameters

$$Im^{out}(x, y) = \sum_{i=-d}^d \sum_{j=-d}^d \sum_{c=1}^C (\textcolor{red}{K(i, j, c)} Im^{in}(x + i, y + j, c) + \textcolor{red}{b_t})$$

- Kernel **parameters**
- $((2d + 1)^2 * C + 1) * T$

# Receptive field

# Convolution

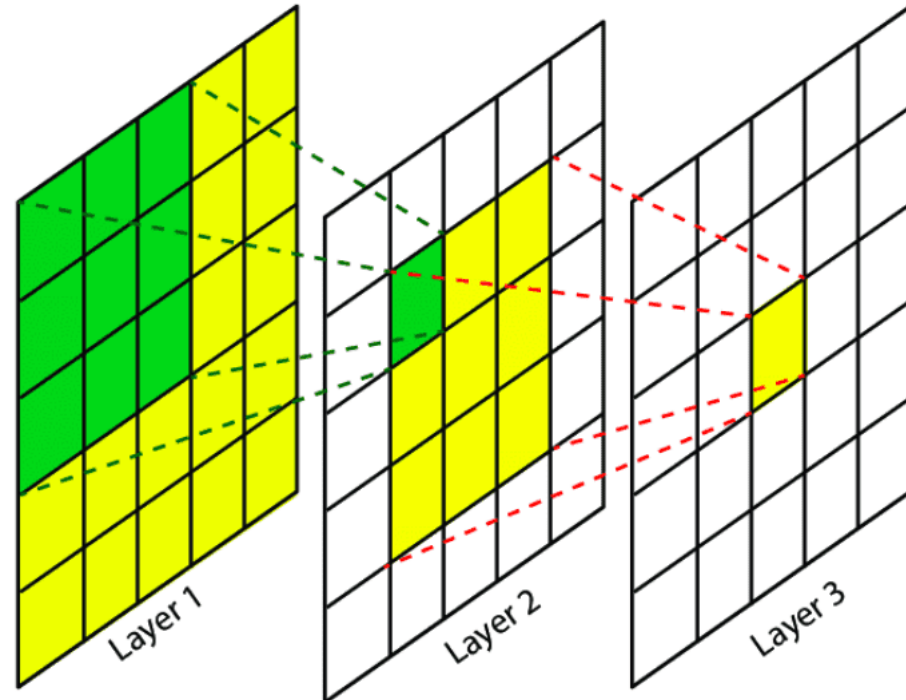


# Receptive field

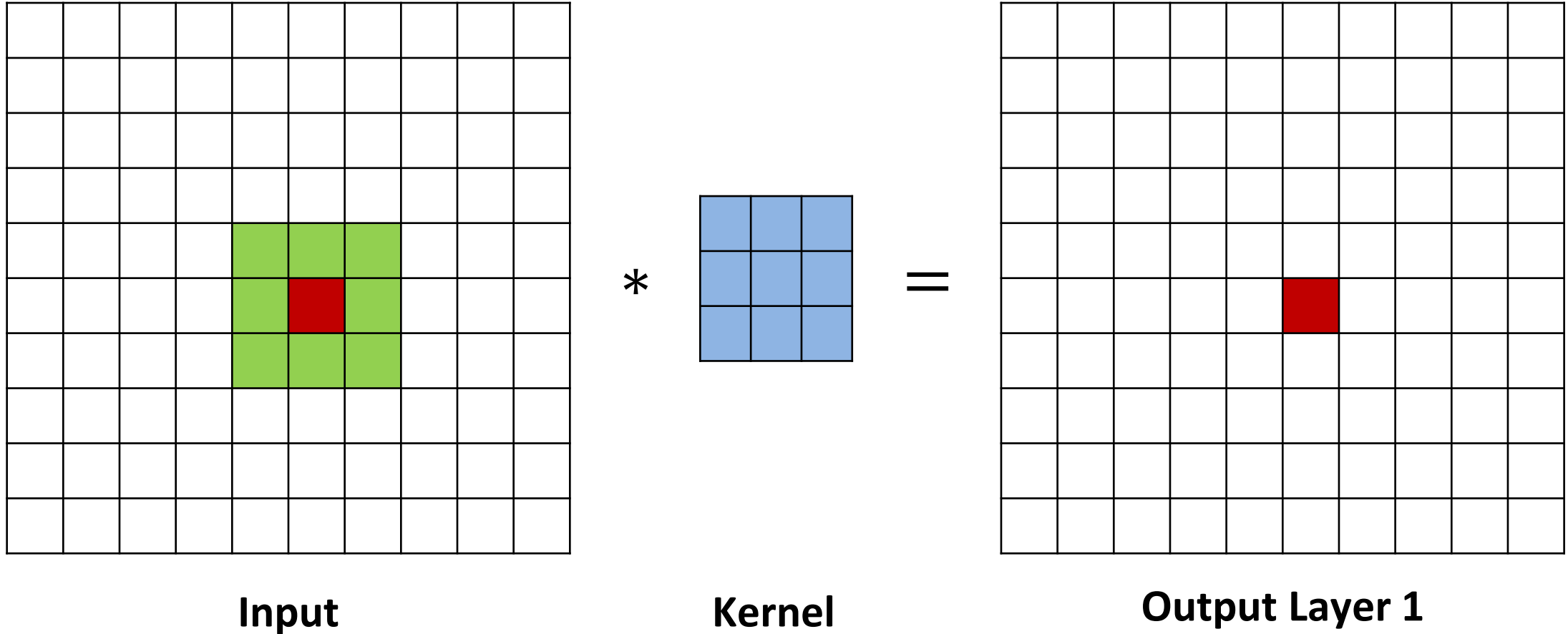
- Consider a pixel in the output image
- What part of the input image influences the value in this output pixel?

## Receptive Field (RF)

the size of the region in the input that produces the feature



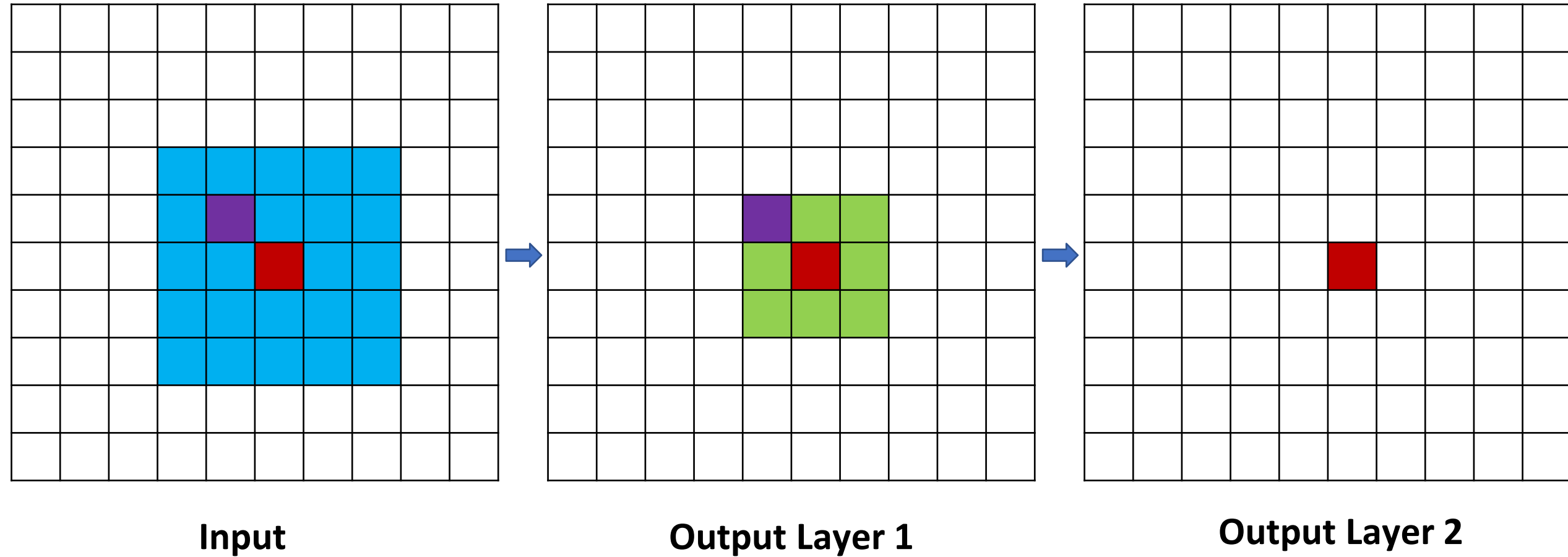
# Receptive field



*Receptive field: 3 x 3*



# Receptive field



*Receptive field: 5 x 5*

# Receptive field

Receptive field for 3 x 3 convolution:

- After 1 convolutional layer: 3 x 3
- After 2 convolutional layers: 5 x 5
- After 3 convolutional layers: 7 x 7

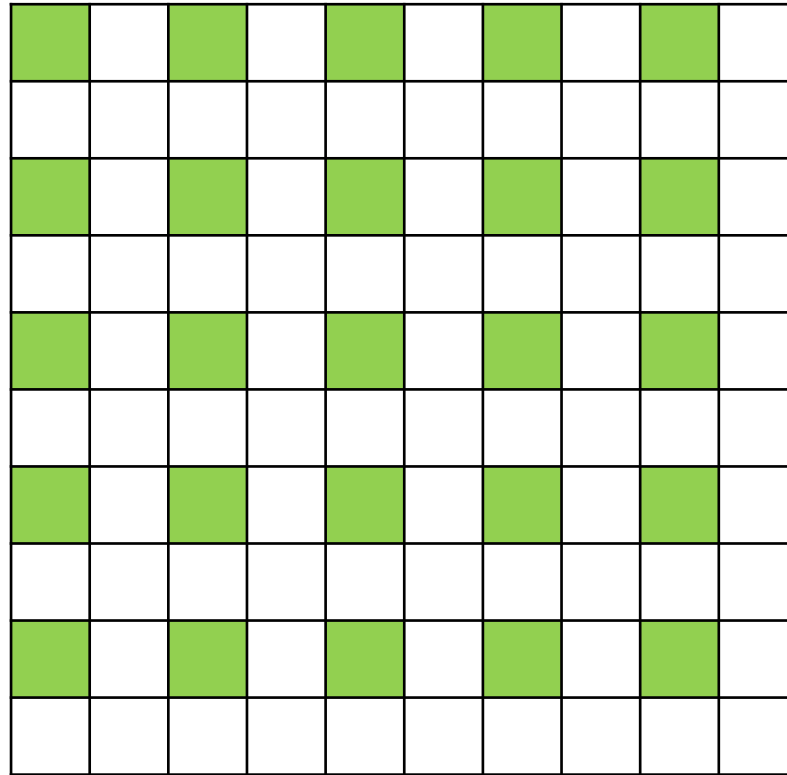
We need a lot of layers if the image size is 512 x 512!

# Strides

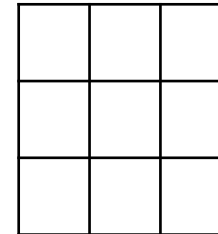
## Stride:

the number of pixels by which we move the filter across the input image

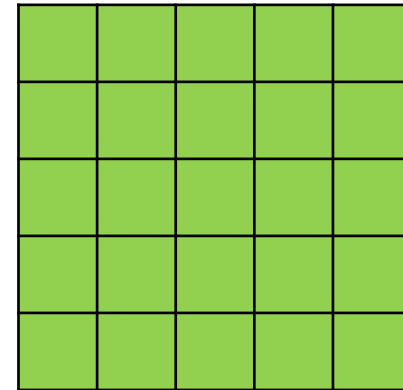
$$s = 2$$



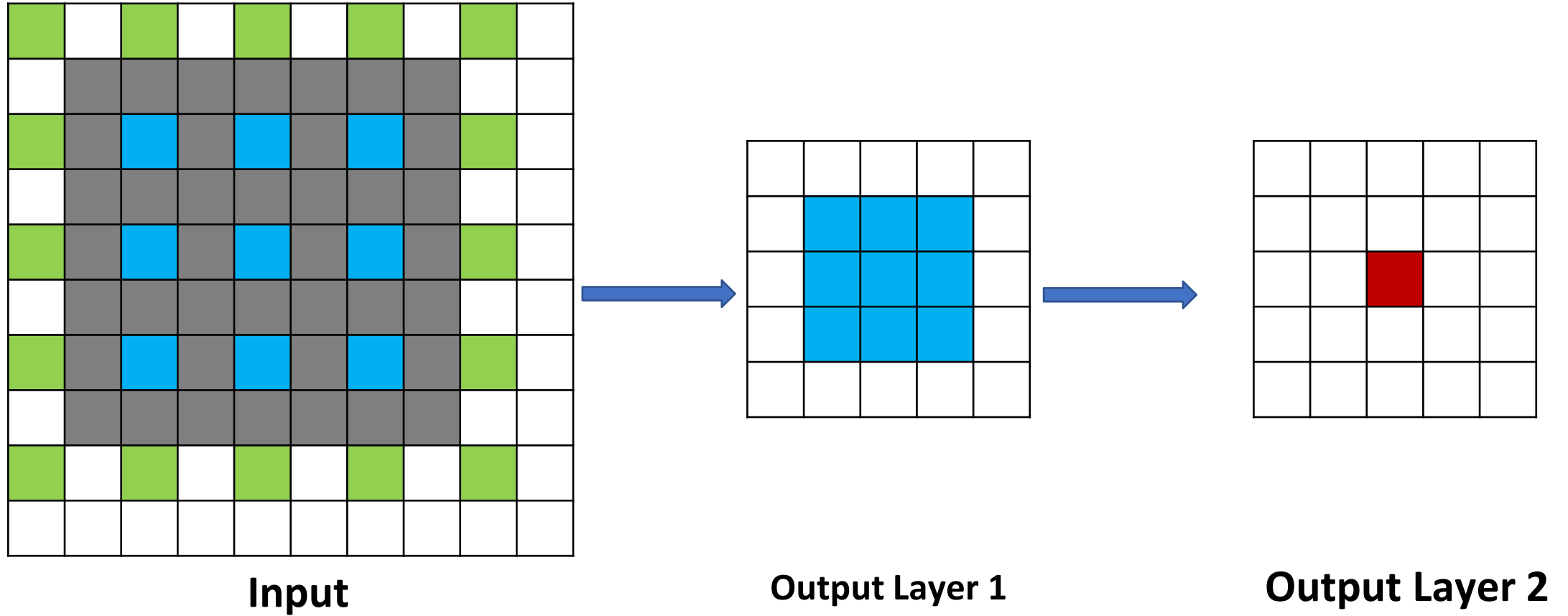
\*



=



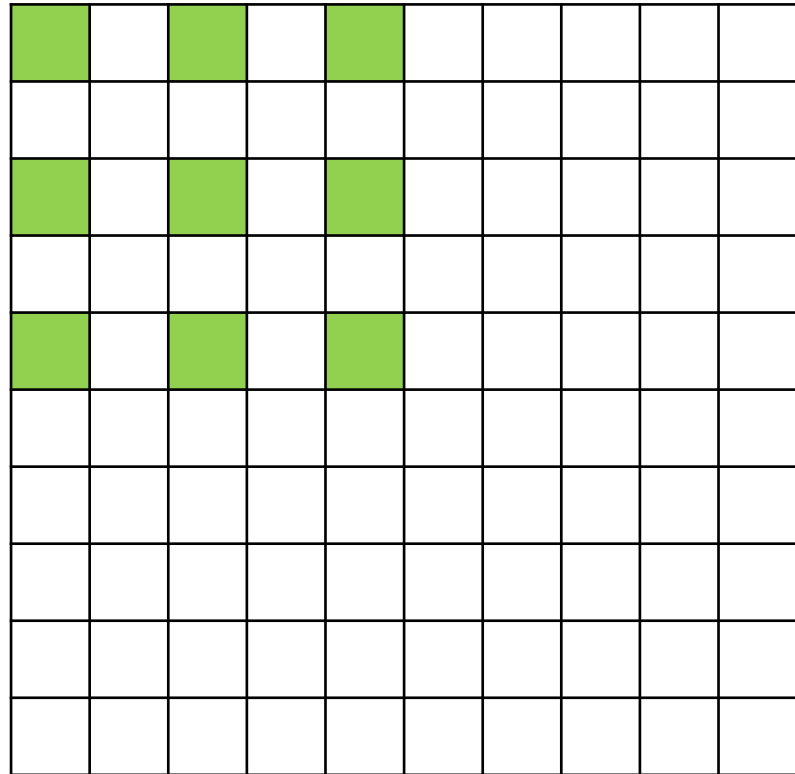
# Strides



*Receptive field: 7 x 7*

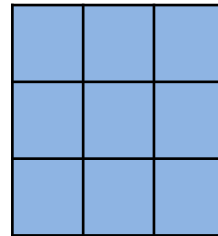
# Dilated convolutions

$l = 2$



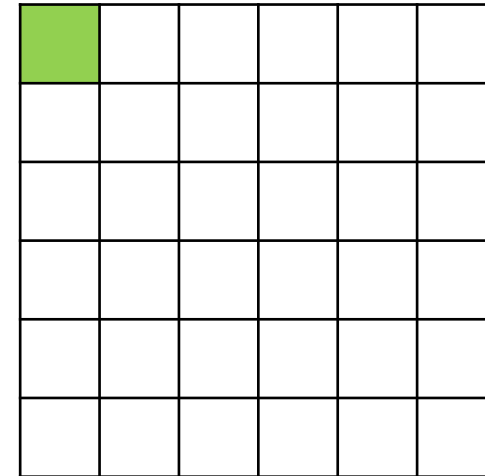
Input

\*



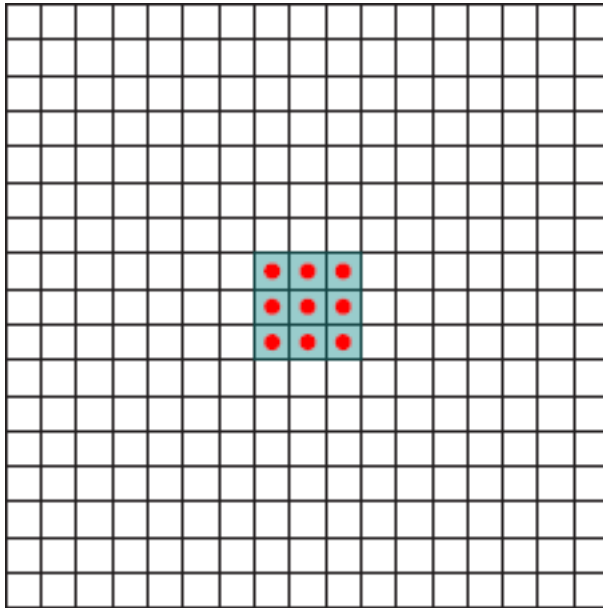
Kernel

=

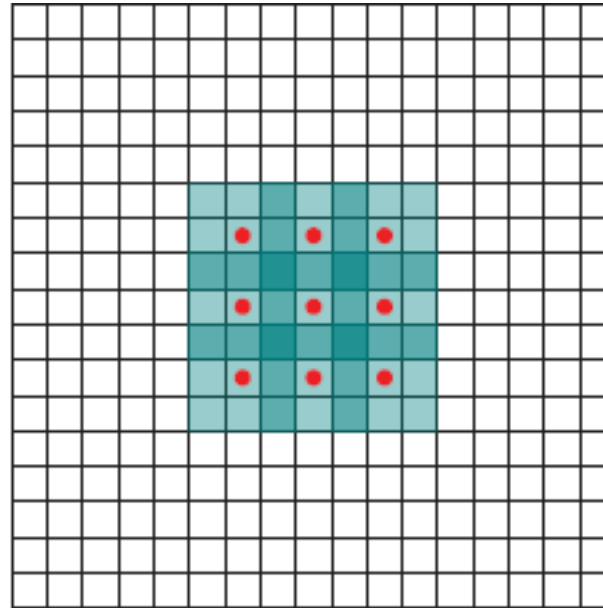


Output Layer 1

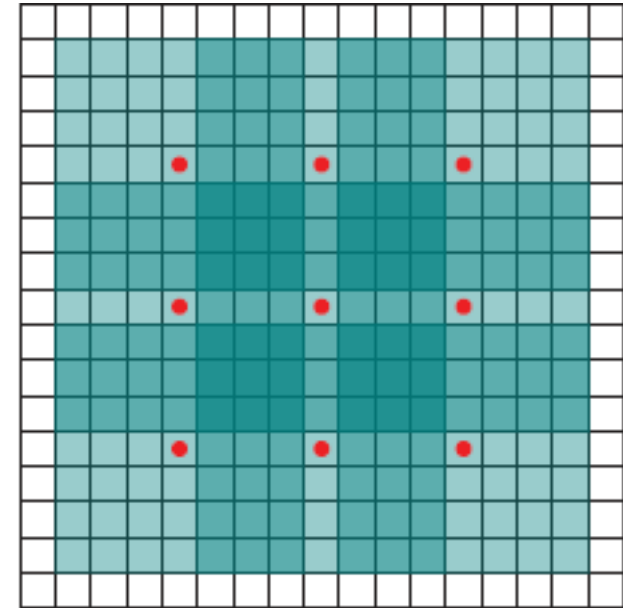
# Dilated convolutions



$l = 1$

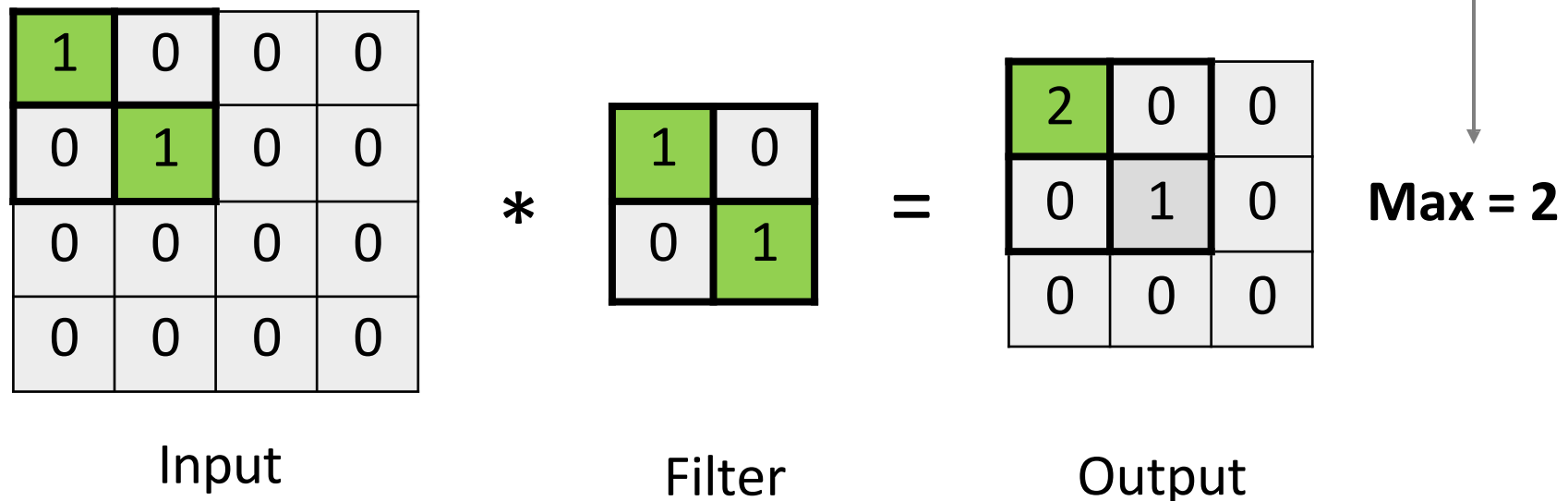
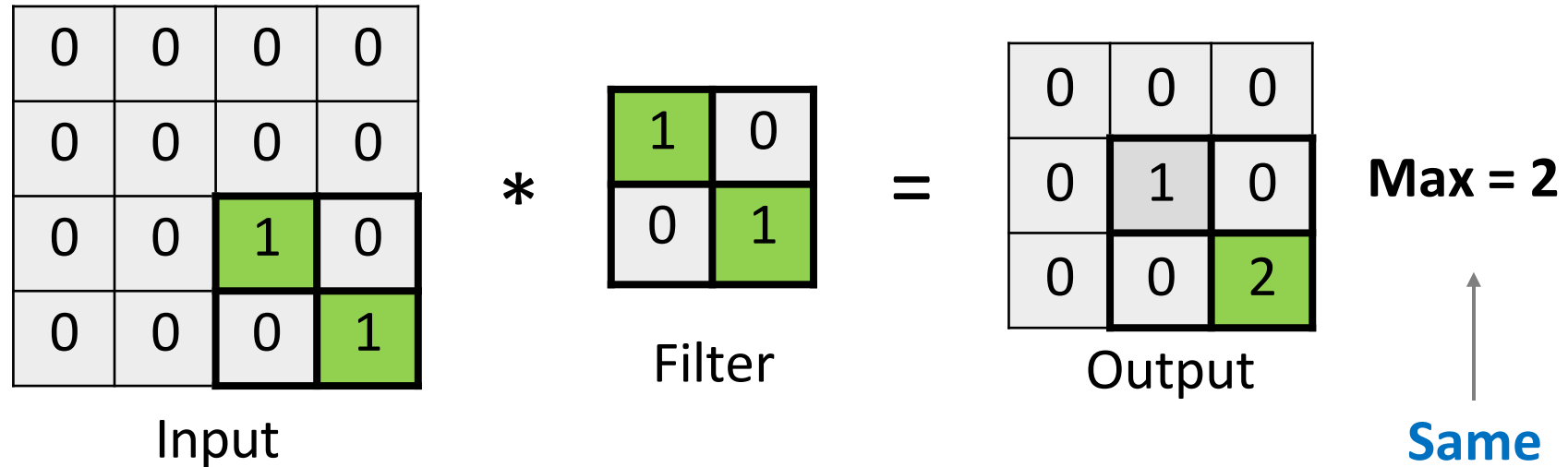


$l = 2$



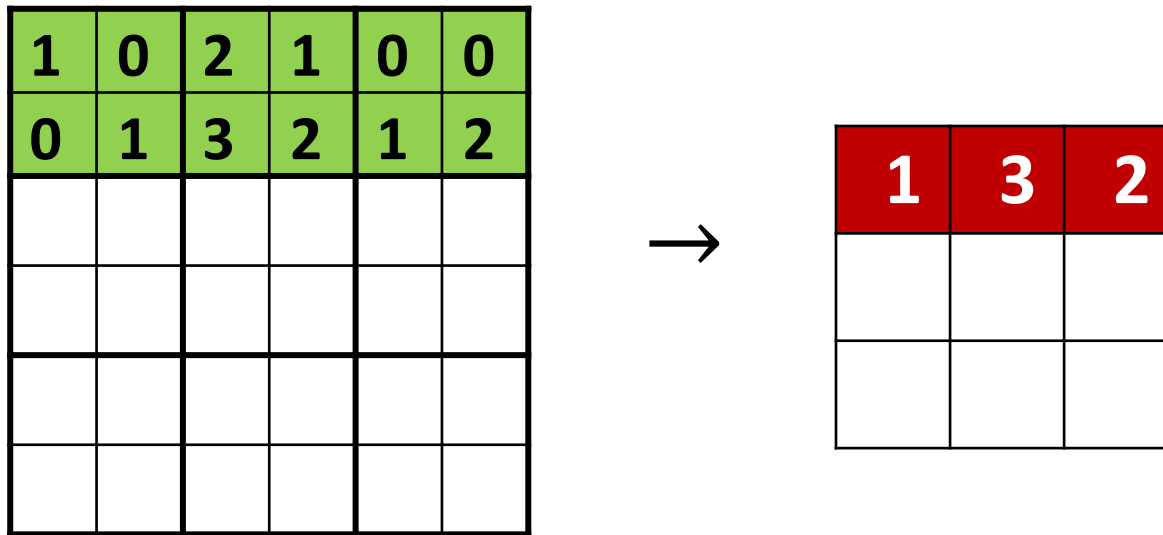
$l = 4$

# The maximum of convolution is invariant to shifts



# Pooling

- Max-pooling with kernel 2x2





# Pooling

- Splits the image into  $n \times m$  sections applying some function (usually max)
- Significantly reduces the size of the image (which means it increases the receptive field in the subsequent layers)
- Has no parameters

1	0	2	1	0	0
0	1	3	2	1	2



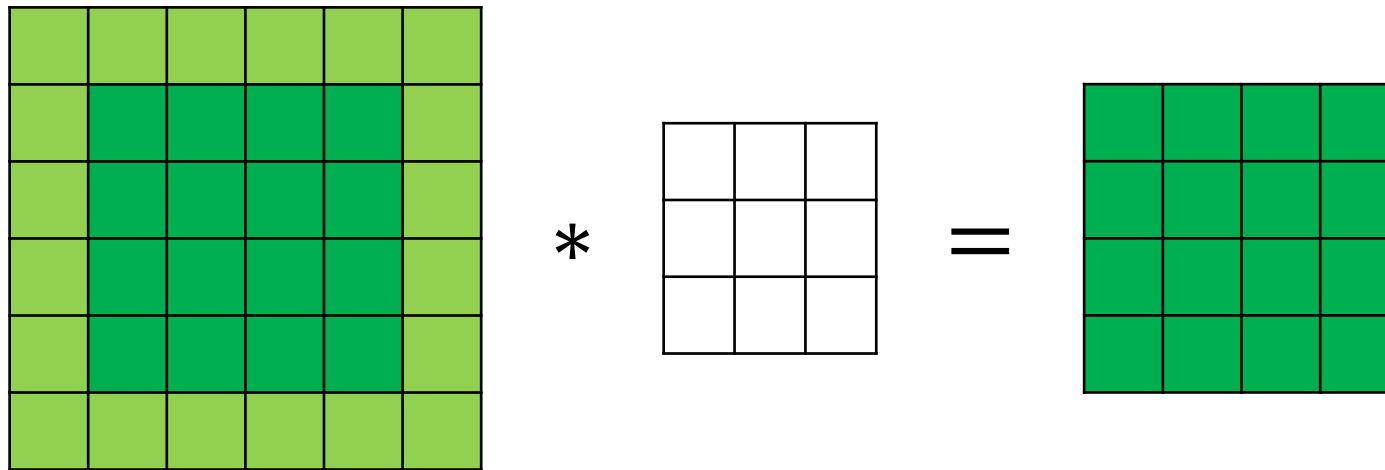
1	3	2

# Why we need to know all this?

- It is important to ensure that the last convolutional layer has a receptive field relative to the size of input image

# Convolution

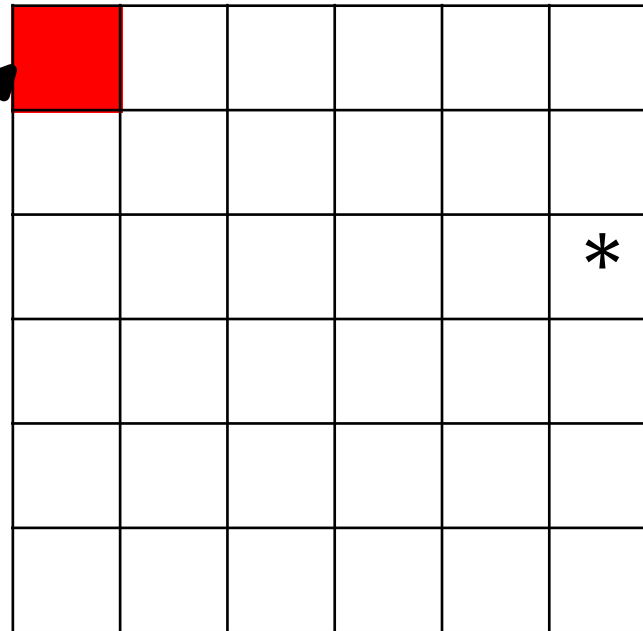
- If you apply convolution, the output image will be smaller than the input



# Valid padding (= No padding)

- When applying convolution to image, the pixels at the edges have less impact on the output

We will not see that the filter has a good response when placing the center at this pixel



0	0	1
0	0	1
1	1	1

# Zero padding

0	0	0	0	0	0	0	0
0							0
0							0
0							0
0							0
0							0
0							0
0	0	0	0	0	0	0	0

\*

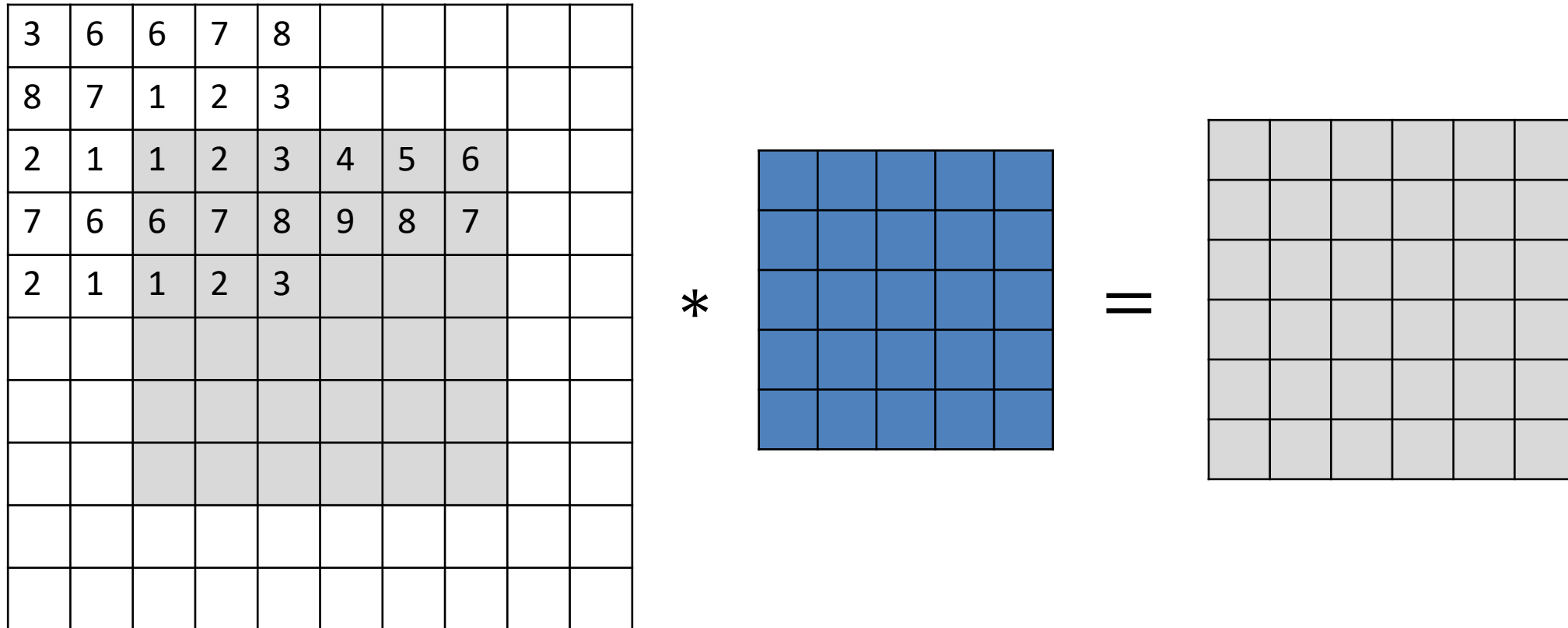

=


# Zero padding

- We add zeros along the boundaries so that the convolution calculated after this in valid mode gives an image of the same size as the original one
- There is a risk that the model will learn to understand where the edges are in the image - we may lose invariance

# Reflection padding

- Pad input data symmetrically with a reflection of its own values



# Reflection padding

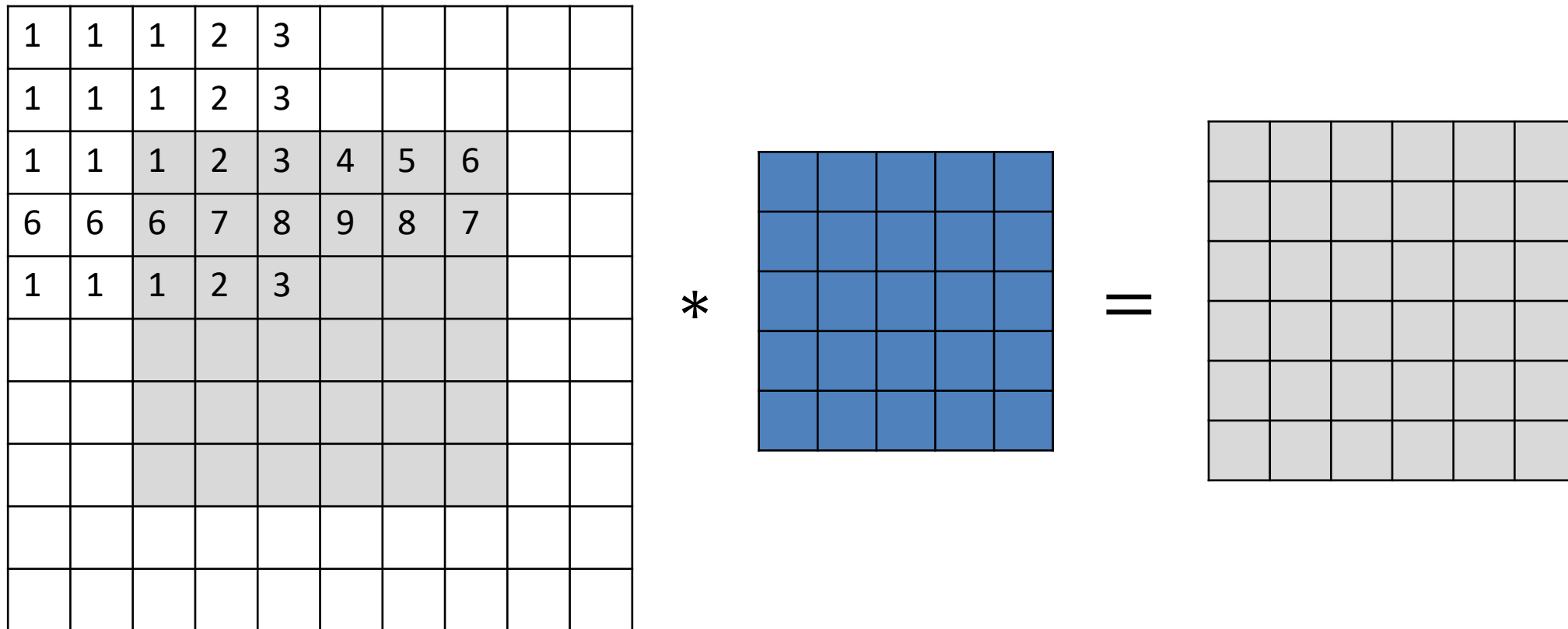
- Can't easily find image edges
- But now the model can begin to find specular reflections and select filters for them

[illegible]



# Replication padding

- Pad with replicated the values at the boundary



# Replication padding

- The pixel on the border is equal to the nearest pixel from the image
- The model can still adjust to the patterns that arise from such padding

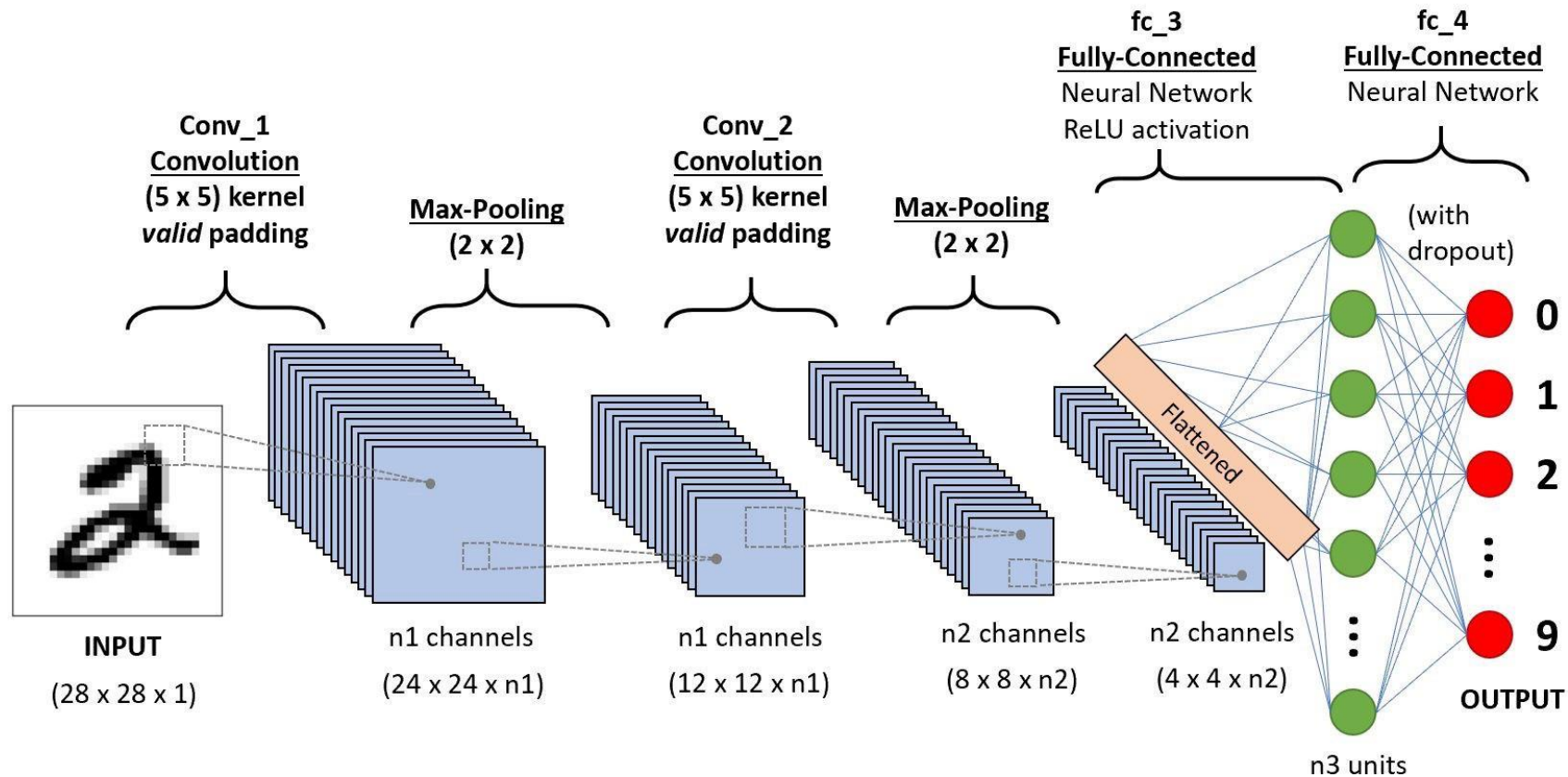
1	1	1	2	3					
1	1	1	2	3					
1	1	1	2	3	4	5	6		
6	6	6	7	8	9	8	7		
1	1	1	2	3					

# Summary

- Padding allows you to control the size of the output images
- Padding allows you to take into account objects on the edges
- Different types of padding allow different methods of retraining for edges

# Convolutional Neural Networks

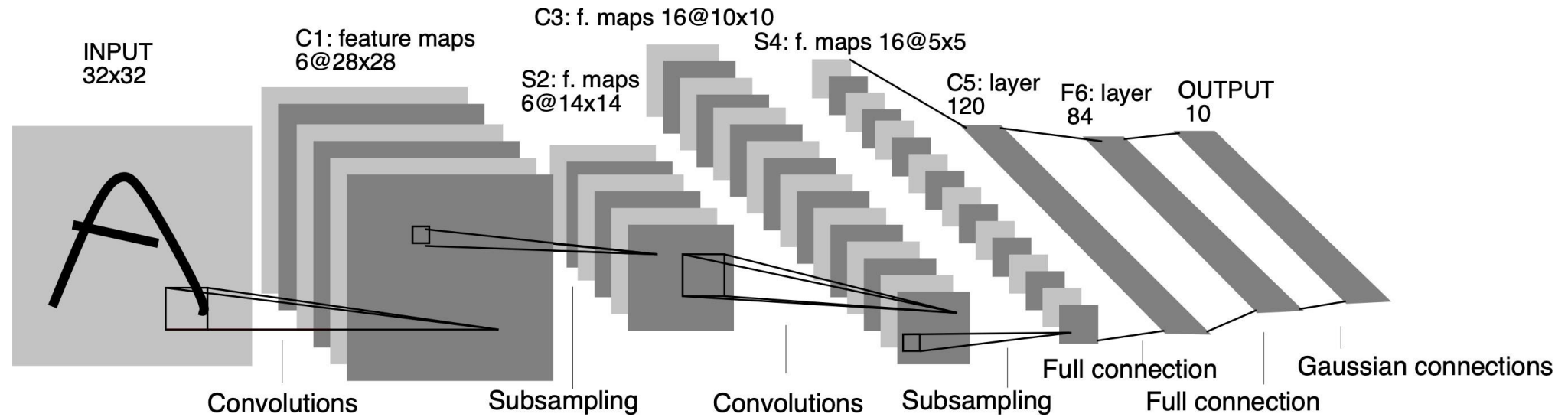
# Architecture 1



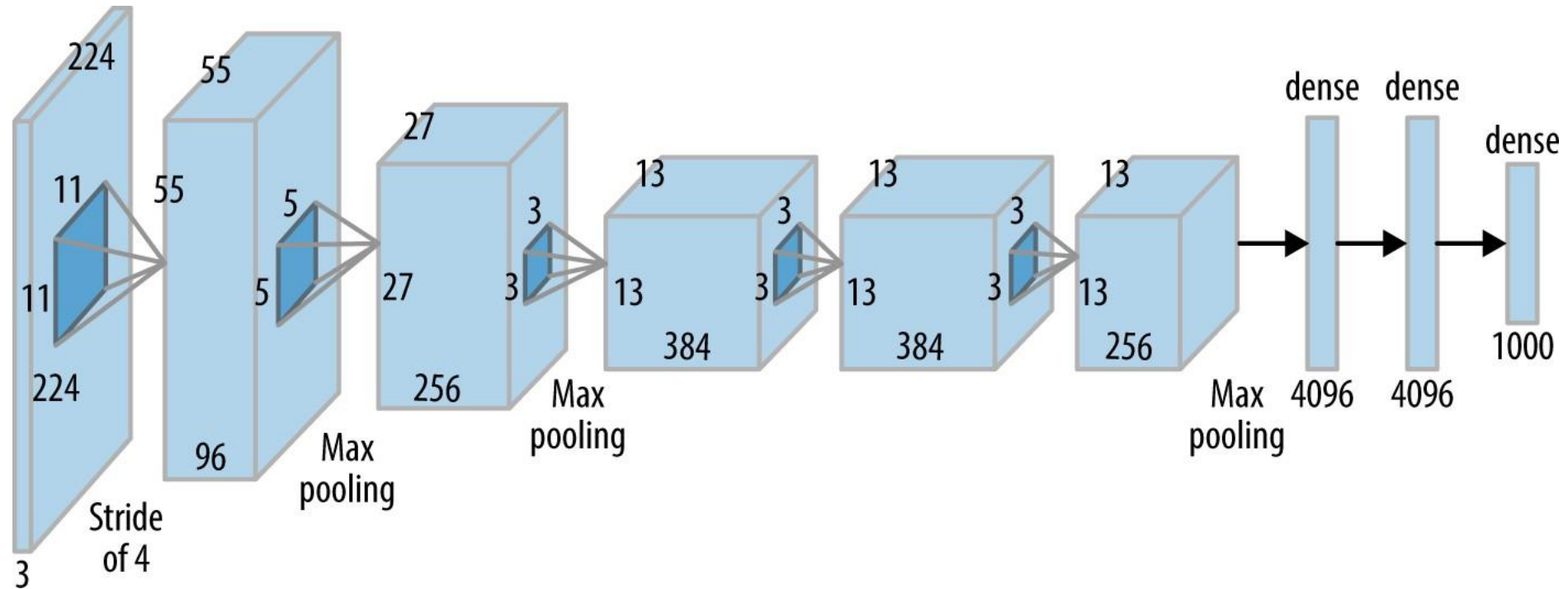
# Architecture 1

- 1) Convolution->linear layer>pooling or convolution->non-linear layer
- 2) Flatten the output
- 3) Fully-connected layer

# Architecture 2: LeNet



# Architecture 3: AlexNet



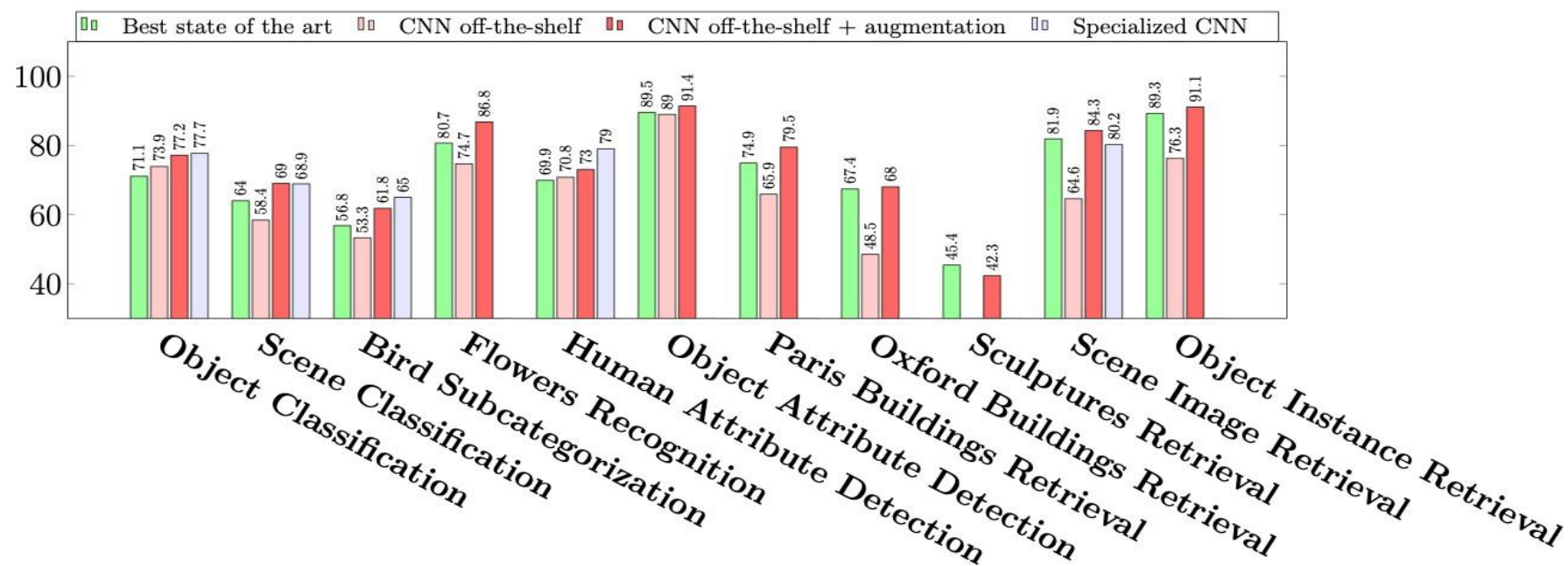
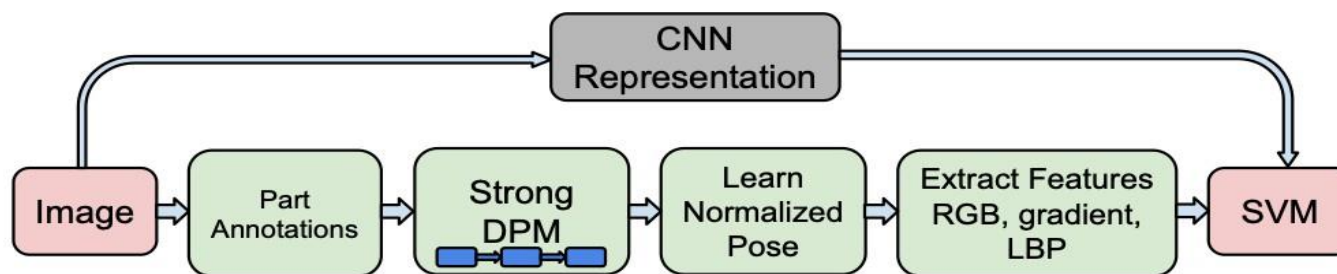


# Visualizing Convolutional Networks

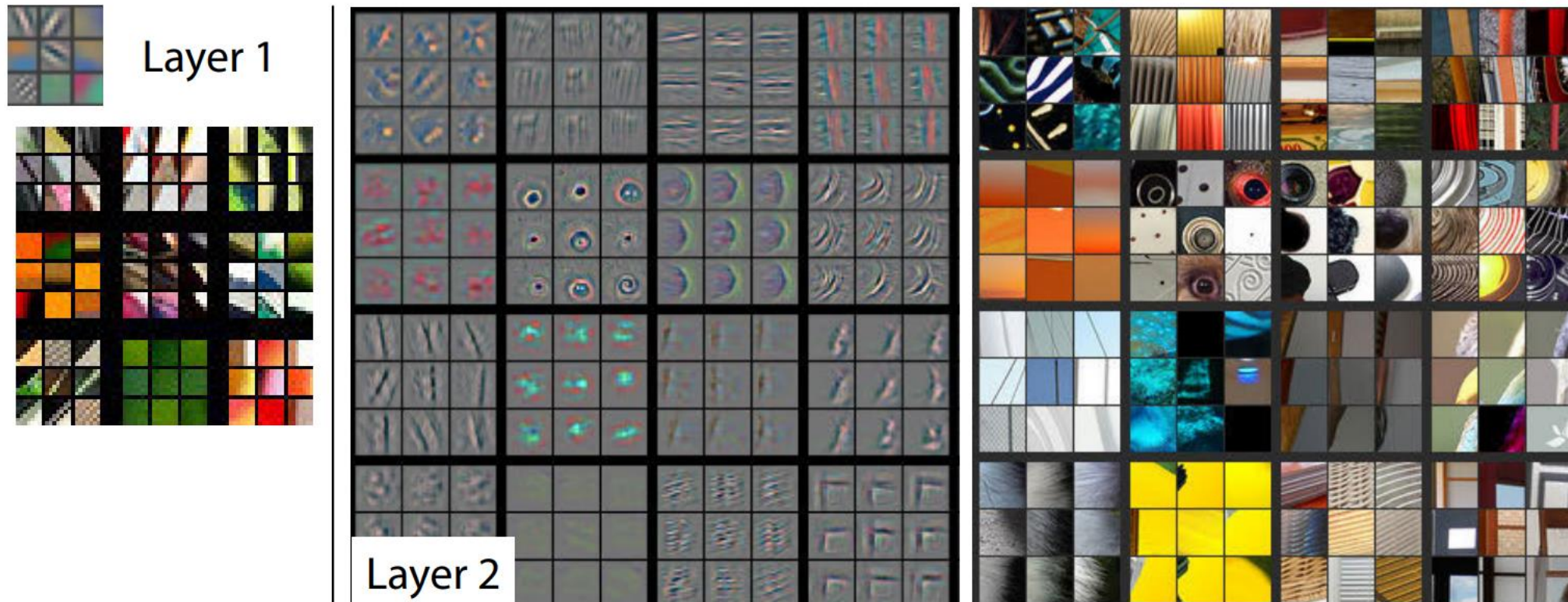
# Image representation (embedding) from the last layers

- Important observation: the output of final layers serve as good feature representations of images
- For instance, they can be used in tasks like searching for similar images

# Last layer embeddings

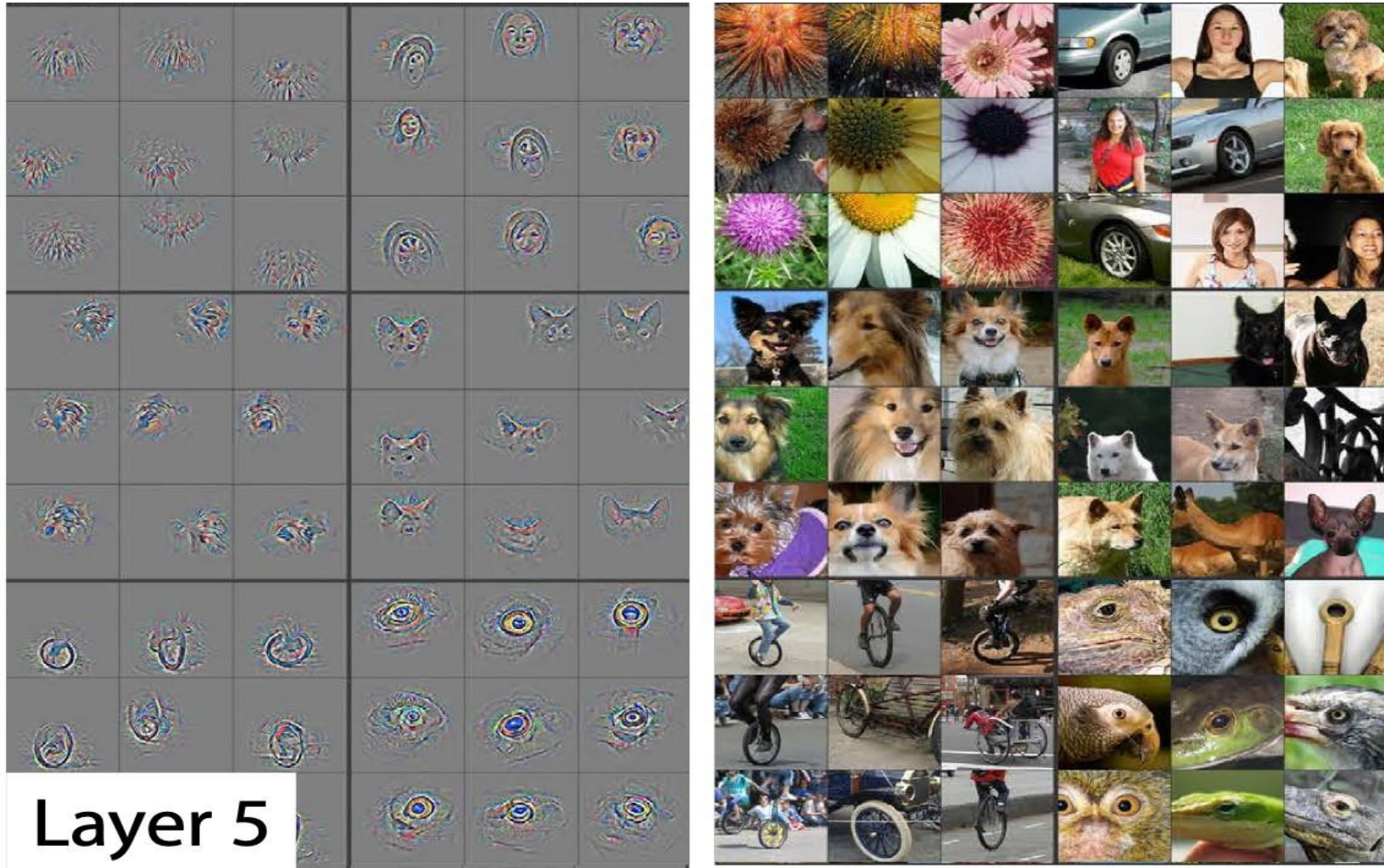


# Visualizing Convolutional Networks





# Visualizing Convolutional Networks



Layer 5